

On Fast Convergence of Proximal Algorithms for SQRT-Lasso Optimization: DonâĂŹt Worry About its Nonsmooth Loss Function



Xingguo Li^{\$}, Haoming Jiang[†], Jarvis Haupt^{*}, Raman Arora^{*}, Han Liu^{\$}, Mingyi Hong^{*}, Tuo Zhao[†] ^oPrinceton University *University of Minnesota *Johns Hopkins University ^oNorthwestern University [†]Georgia Tech

Geometric Properties of SQRT-Lasso

Lemma 1. Given sub-Gaussian assumption of X and ϵ , and $||\theta^*||_0 = s^*$, then with high probability: (i) for $\forall \lambda \geq C_1 \sqrt{\frac{\log d}{n}}$, we have

$$\lambda \ge \frac{C_1}{4} ||\nabla \mathcal{L}(\theta^*)||_{\infty},$$

which is sufficiently large to yield sufficiently sparse solutions;

(ii) for $n \ge C_2 s^* \log d$, $\exists \mathcal{B}_r = \{\theta \in \mathbb{R}^d : ||\theta - \theta^*||_2^2 \le r\}$ for some constant r s.t. for $\forall v, w \in \mathcal{B}_r$ satisfying $||v - w||_0 \leq s^* + 2\widetilde{s}$, \mathcal{L} is locally restricted strongly convex (LRSC), smooth (LRSS), and Hessian smooth (LRHS), if $\exists \rho_s^-, \rho_s^+, L_s \in (0, \infty)$ i.e.

$$LRSC:\mathcal{L}(v) - \mathcal{L}(w) - \nabla \mathcal{L}(w)^{\top}(v - w) \ge \frac{\rho_{s}^{-}}{2} ||v - w||_{2}^{2},$$

$$LRSS:\mathcal{L}(v) - \mathcal{L}(w) - \nabla \mathcal{L}(w)^{\top}(v - w) \le \frac{\rho_{s}^{+}}{2} ||v - w||_{2}^{2},$$

$$LRHS: \min_{\substack{||u||_{0} \le s, \\ ||u||_{2} = 1}} u^{\top} (\nabla^{2} \mathcal{L}(v) - \nabla^{2} \mathcal{L}(w)) u \le L_{s} ||v - w||_{2}^{2}$$

with

$$\rho_{s^*+2\widetilde{s}}^+ \leq \frac{C_3}{\sigma}, \ \rho_{s^*+2\widetilde{s}}^- \geq \frac{C_4}{\sigma} \ \text{and} \ L_{s^*+2\widetilde{s}} \leq \frac{C_5}{\sigma}$$

where $C_1, \ldots, C_5 \in \mathbb{R}^+$ are generic constants.

Computational Theory

Denote

$$\kappa_s = \frac{\rho_s^+}{\rho_s^-}, \quad \mathcal{S}^* = \{j \mid \theta_j^* \neq 0\}, \quad \overline{\mathcal{S}}^* = \{j \mid \theta_j^* = 0\}, \\ \mathcal{B}_r^{s^* + \widetilde{s}} = \mathcal{B}_r \cap \{\theta \in \mathbb{R}^d : ||\theta - \theta^*||_0 \le s^* + \widetilde{s}\}.$$

Theorem 2 (Local Convergence). Suppose X and n satisfy Lemma 1. Given λ and $\theta^{(0)}$ such that $\lambda \geq \frac{C_1}{4} || \nabla \mathcal{L}(\theta^*) ||_{\infty}$, $|| heta^{(0)} - heta^*||_2^2 \leq s^* \left(8\lambda/
ho_{s^*+\widetilde{s}}^ight)^2$ and $heta^{(0)} \in \mathcal{B}_r^{s^*+\widetilde{s}}$, we have sufficiently sparse solutions for all iterations, i.e.,

 $||[\theta^{(t)}]_{\overline{\mathcal{S}}^*}||_0 \le \widetilde{s}.$

Moreover, given $\varepsilon > 0$, we need at most

$$T = \mathcal{O}\left(\kappa_{s^*+2\widetilde{s}}\log\left(\frac{\kappa_{s^*+2\widetilde{s}}^3 s^* \lambda^2}{\varepsilon^2}\right)\right) \qquad \text{Prox-Grad}$$
$$T = \mathcal{O}\left(\log\log\left(\frac{3\rho_{s^*+2\widetilde{s}}^+}{\varepsilon}\right)\right) \qquad \text{Prox-Newton}$$

iterations to guarantee that the output solution θ satisfies

$$\begin{split} ||\widehat{\theta} - \overline{\theta}||_{2}^{2} &= \mathcal{O}\left(\left(1 - \frac{1}{8\kappa_{s^{*}+2\widetilde{s}}}\right)^{T} \varepsilon \lambda s^{*}\right) \qquad \text{Prox-Grad} \\ ||\widehat{\theta} - \overline{\theta}||_{2}^{2} &= \mathcal{O}\left(\left(\frac{L_{s^{*}+2\widetilde{s}}}{2\rho_{s^{*}+2\widetilde{s}}}\right)^{2^{T}} \varepsilon \lambda s^{*}\right) \qquad \text{Prox-Newton,} \end{split}$$

where $\overline{\theta}$ is the unique sparse global optimum to (1) with $||[\overline{\theta}]_{\overline{S}^*}||_0 \leq |\overline{S}^*||_0 \leq |\overline{B}|_{\overline{S}^*}||_0 < |\overline{B}|_0 < |\overline{B}|_0$

Statistical Theory

Theorem 3. Suppose X, and n satisfy conditions in Lemma 1. Given $\lambda = C_1 \sqrt{\log d/n}$, if the output solution $\widehat{ heta}$ obtained from Algorithm 1 and 2 satisfying, $\omega_{\lambda}(\widehat{ heta}) \leq \varepsilon = \mathcal{O}\left(rac{\sigma s^* \log d}{n}
ight)$, then we have:

In Algorithm 1,2, $\hat{\theta}$ achieves the minimax optimal rate of convergence in parameter estimation.

Input

Initial For: *I*

End F Retur



$$\begin{split} ||\widehat{\theta} - \theta^*||_2 &= \mathcal{O}_P\left(\sigma\sqrt{\frac{s^*\log d}{n}}\right) \quad \text{and} \\ ||\widehat{\theta} - \theta^*||_1 &= \mathcal{O}_P\left(\sigma s^*\sqrt{\frac{\log d}{n}}\right). \end{split}$$

Moreover, we have $|\widehat{\sigma} - \sigma| = \mathcal{O}_P\left(rac{\sigma s^* \log d}{n}
ight)$, where $\widehat{\sigma} = rac{||y - X\widehat{ heta}||_2}{\sqrt{n}}$.

Boosting via Pathwise Optimization Scheme

Algorithm 3 The pathwise optimization scheme

Theorem 4. Suppose the design matrix X is sub-Gaussian, and $\lambda_{[N]} = C_1 \sqrt{\log d/n}$. For $n \geq C_2 s^* \log d$ and $\eta_{\lambda} \in (\frac{5}{6}, 1)$, the following results hold:

(I) There exists an $N_1 < N$ such that

$$r > s^* \left(8\lambda_{N_1} / \rho_{s^* + \widetilde{s}}^- \right)^2;$$

(II) For any $K \in [N_1 + 1, .., N]$, we have

$$|\theta_{[K]}^{(0)} - \theta^*||_2^2 \le s^* \left(8\lambda_{[K]}/\rho_{s^*+\tilde{s}}^-\right)^2, \theta_{[K]}^{(0)} \in \mathcal{B}_r^{s^*+\tilde{s}};$$

(III) Theorems 2 hold for all λ_K 's, where $K \in [N_1 + 1, .., N]$.

Experiments

Datasets:

- "DrivFace" (n = 606, d = 6400, [1])
- simulated datasets: (n = 1000, d = 1k, 5k, 10k)

the stopping criterion ε_N .





- "Greenhouse" (n = 2921, d = 5232, [2])
- Tab. 1 Computational performance of Prox-GD on synthetic data under different choices of variance σ , the number of stages N, and

N		Minimal		
1 V	10^{-4}	10^{-5}	10^{-6}	MSE
1	0.372	0.372	0.365	
10	0.275	0.276	0.280	0.013
30	0.336	0.345	0.351	
1	0.235	0.248	0.262	
10	0.104	0.103	0.109	1.183
30	0.217	0.222	0.220	

SQRT-Lasso						
rox-GD	Newton	ADMM	ScalReg	CD	Alt.Min	PISTA
5.812	1.708	1027	3181	14.31	99.81	5.113
0.421	0.426	18.88	124.0	3.138	17.69	0.414

Tab. 3 Extension to calibrated multivariate regression.

Λ_N Prox-GDNewtonADMMCD $\sqrt{\log d/n}$ 0.29640.032014.832.410 $2\sqrt{\log d/n}$ 0.17250.02132.2312.227 $4\sqrt{\log d/n}$ 0.04780.01121.8681.366		Synthetic ($\sigma = 1$)					
$\sqrt{\log d/n}$ 0.2964 0.0320 14.83 2.410 $2\sqrt{\log d/n}$ 0.1725 0.0213 2.231 2.227 $4\sqrt{\log d/n}$ 0.0478 0.0112 1.868 1.366	λ_{N}	Prox-GD	Newton	ADMM	CD		
$2\sqrt{\log d/n}$ 0.1725 0.0213 2.231 2.227 $4\sqrt{\log d/n}$ 0.0478 0.0112 1.868 1.366	$\sqrt{\log d/n}$	0.2964	0.0320	14.83	2.410		
$4\sqrt{\log d/n}$ 0.0478 0.0112 1.868 1.366	$2\sqrt{\log d/n}$	0.1725	0.0213	2.231	2.227		
\mathbf{v} \mathbf{C} /	$4\sqrt{\log d/n}$	0.0478	0.0112	1.868	1.366		
λ_N DrivFace	λ_N	DrivFace					
$\sqrt{\log d/n}$ 9.562 0.2186 158.9 12.77	$\sqrt{\log d/n}$	9.562	0.2186	158.9	12.77		
$2\sqrt{\log d/n}$ 8.688 0.1603 129.4 20.42	$2\sqrt{\log d/n}$	8.688	0.1603	129.4	20.42		
$4\sqrt{\log d/n}$ 1.824 0.0924 94.37 19.17	$4\sqrt{\log d/n}$	1.824	0.0924	94.37	19.17		

[1] K. Diaz-Chito, A. Hernández-Sabaté, and A. M. López. A reduced feature set for driver head pose estimation. Appl. Soft Comput., 45(C):98-107,

[2] D. D. Lucas, C. Yver Kwok, P. Cameron-Smith, H. Graven, D. Bergmann T. P. Guilderson, R. Weiss, and R. Keeling. Designing optimal greenhouse gas observing networks that consider performance and cost. Geoscientific Instrumentation, Methods and Data Systems, 4(1):121–137, 2015.

